# CSE422

## ARTIFICIAL INTELLIGENCE
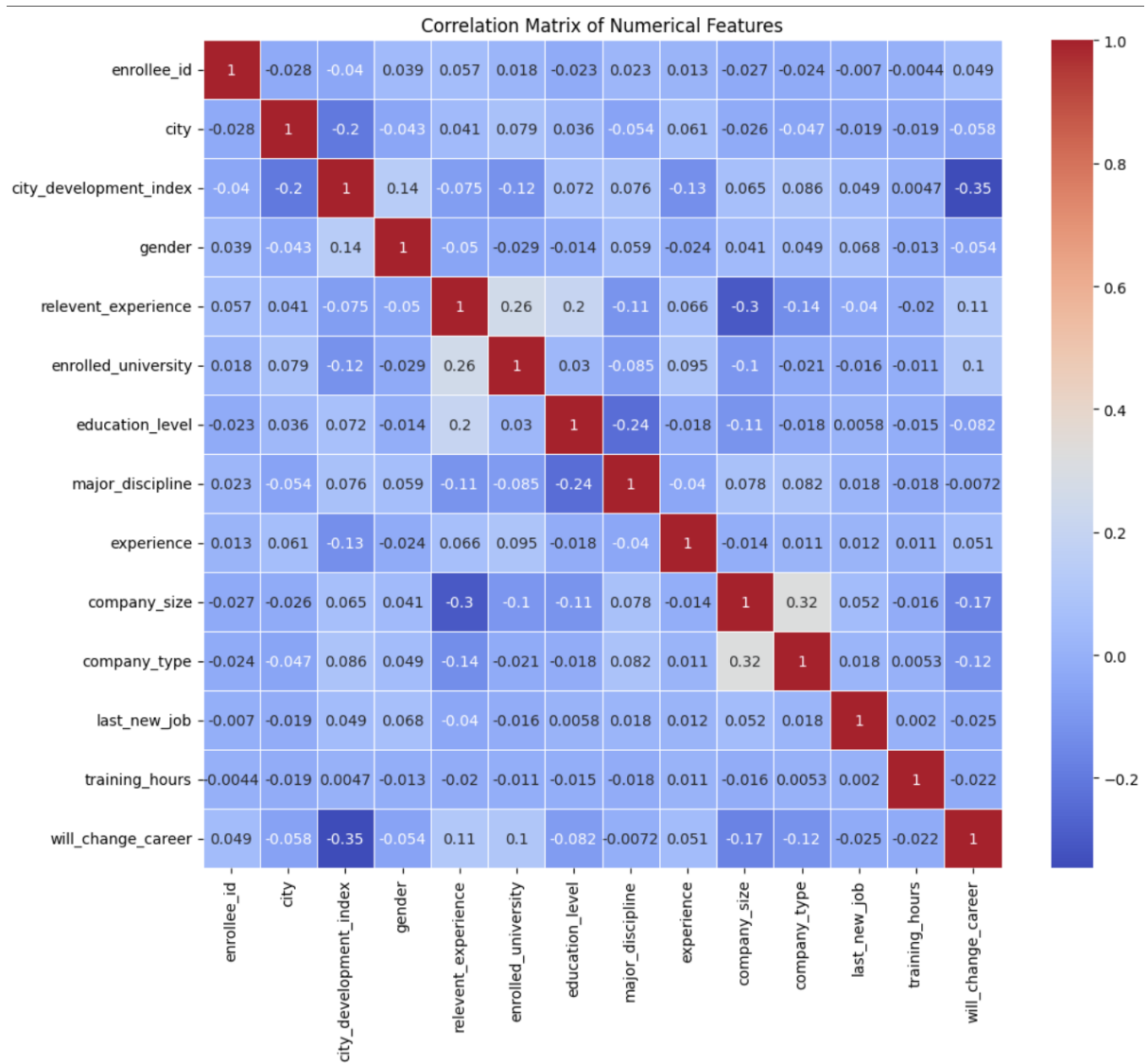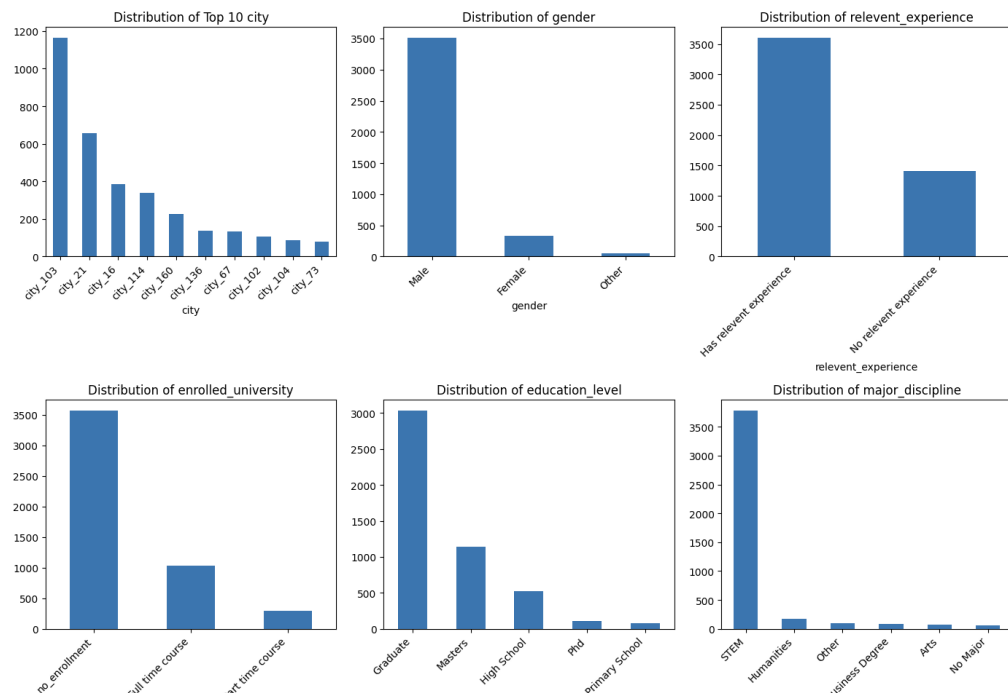
## Table of Content

**1.Introduction:**  In today's rapidly evolving job market, career switching has become increasingly common due to changing interests, economic shifts, or industry demands. Our project focuses on predicting the likelihood of an individual switching careers based on various factors such as enrolled ID, city, city development index, gender, relevant experience, enrolled university, education level, major discipline, experience, company size, company type, last new job, training hours, will change career. The aim is to analyse personal and professional attributes to assess the risk or probability of a career change, which can aid both individuals and organizations in making informed decisions. The motivation behind this project is to apply machine learning techniques to understand the patterns and key indicators that lead to a career switch, ultimately helping career counsellors, HR departments, and job seekers proactively address professional development and career planning.

**2. Dataset Description:** This dataset presents a classification problem where the target variable "will change career" consists of discrete class labels indicating whether an individual is likely to switch careers or not. The dataset contains 14 features and about 7000 data points from there few of them are categorical features among these, several features are categorical, such as city , gender, relevant experience, enrolled university, education level, major discipline, company size, company type, and last new job. Additionally, there are numerical and qualitative features like enrolled ID , city development index, training hours, and years of experience. One of the key considerations in working with this dataset is the potential class imbalance in the target variable, which needs to be addressed to ensure robust model performance.
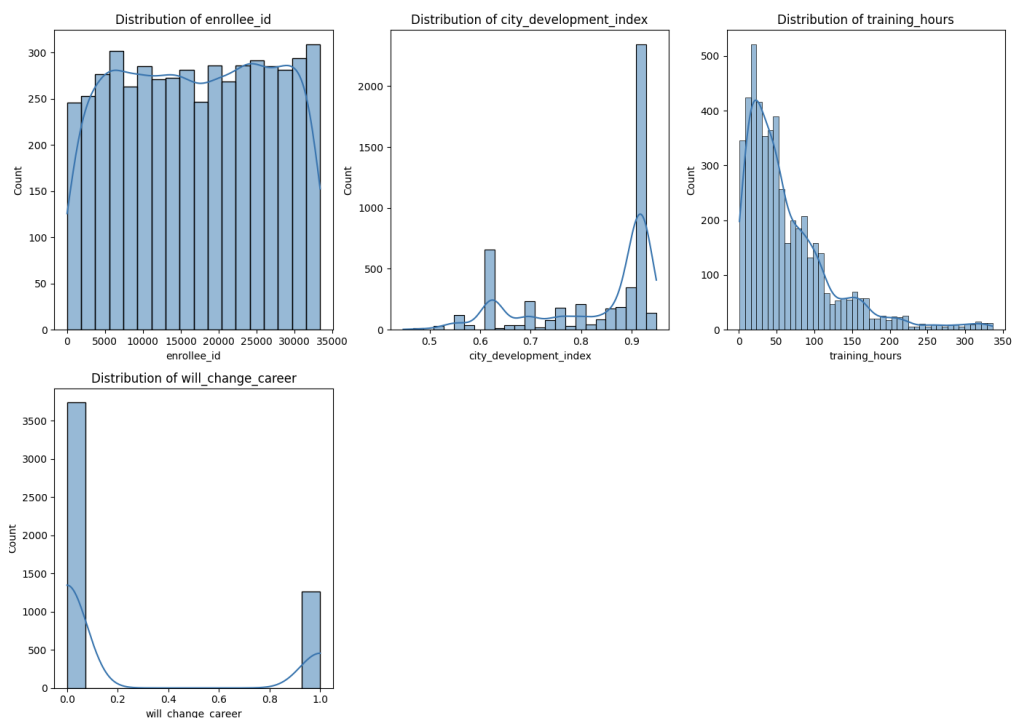
## Correlation:



Correlation Matrix of Numerical Features

# Plot Distribution of Categorical Values



Distribution of Top 10 city

Distribution of gender

Distribution of relevent_experience

Distribution of enrolled_university

Distribution of education_level

Distribution of major_discipline

# Plot Distribution of Numerical Values



Distribution of enrollee_id

Distribution of city_development_index

Distribution of training_hours

Distribution of will_change_career

### 3. Dataset Pre-processing:

It is important to inspect the dataset before training machine learning models to avoid some potential issues. In the dataset there might be some null values. Now from the dataset we have separate the features and target value or label. By runnig this command df.isnull( ) we saw that there are some null values in our dataset. To minimize those null values, for categorical features we replaced the value with the mode (most frequent value) and for numerical features we replaced the value with median. We also noticed that there are some dtype = object values that means we got categorical features in our dataset. Machine learning algorithm such as logistic regression and KNN require numerical input. Therefore, categorical variables must be transformed into numerical formats to ensure compatibility with the models. To solve this problem, we used encoding to convert categorical features into numerical format. This creates new binary columns for each category in a categorical feature.

**Feature Scaling:** We applied feature scaling to ensure that all features have a comparable scale, which is essential for distance-based models like KNN and gradient-based methods like Logistic Regression. We used Standardization to scale the features. This method transforms the data to have a mean of 0 and a standard deviation of 1.

**4. Dataset Splitting:** The dataset was split into training and testing sets. We used an 70-30 split, where 70% of the data was used for training and 30% for testing. This was implemented using the train_test_split function from sklearn. We also used random start and stratify parameters to improve the reliability of the model evaluation.

**5. Model Training and Testing:** Three models were trained and tested: Logistic Regression, K-Nearest Neighbors (KNN), and Decision Tree. We implemented these models to compare the performance of these models on a multiclass classification problem and determine the best-performing one.
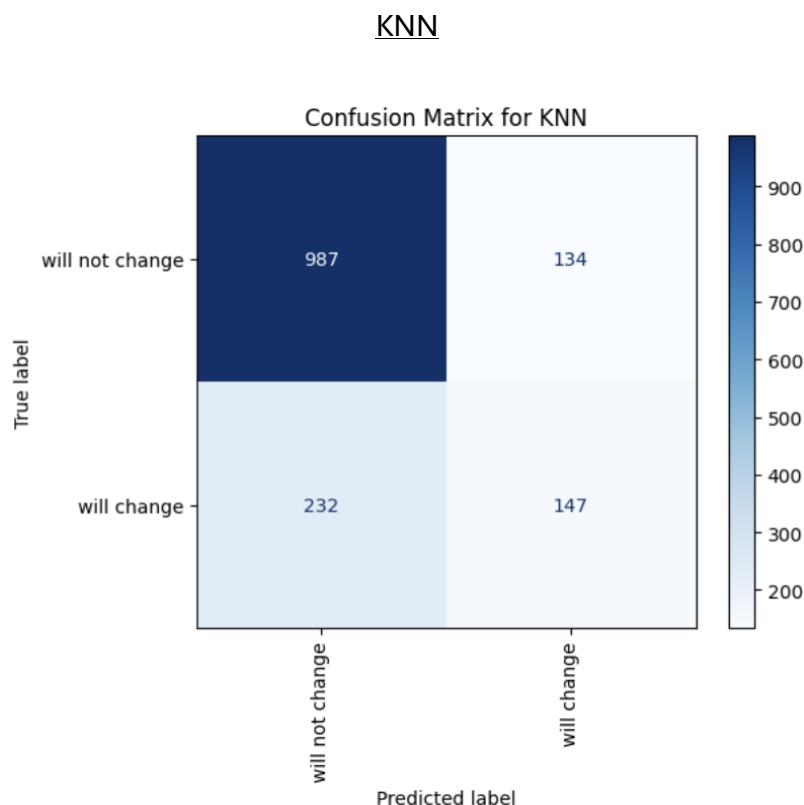
- **K-Nearest Neighbors (KNN):** A distance-based algorithm where predictions are made based on the closest k neighbors in the feature space.
- **Decision Tree:** A tree-based algorithm that splits data based on feature values to make predictions.
- **Logistic Regression:** A linear model for classification problems.

Each model was trained on the training set (X_train, y_train) and tested on the test set (X_test, y_test).
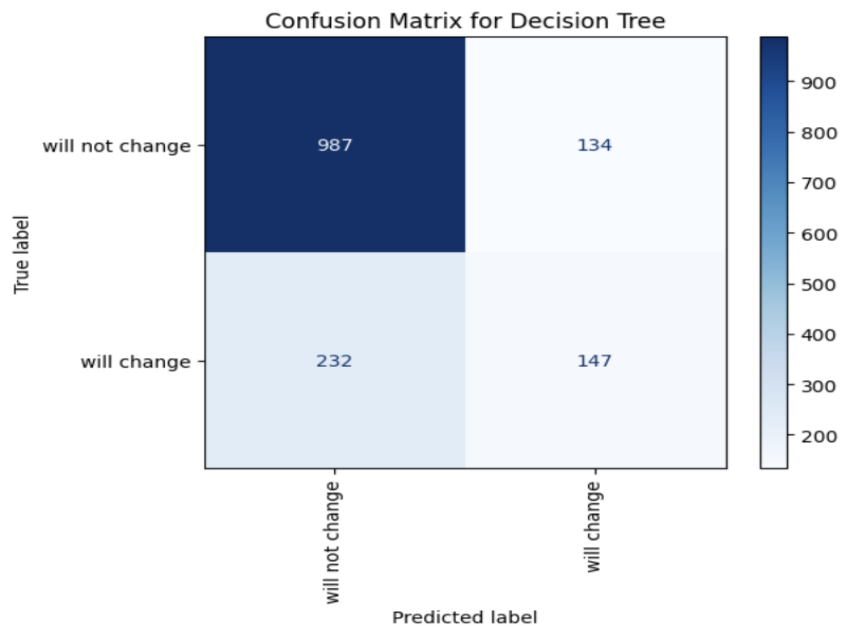
**Comparison:**

| Model Name | Accuracy | Precision | Recall |
|---|---|---|---|
| KNN | 74 | 0.49 | 0.33 |
| Decision Tree | 72 | 0.43 | 0.43 |
| Logistic Regression | 77 | 0.59 | 0.26 |
| Neural Network | 78 | 0.61 | 0.35 |

# Confusion Matrix:

KNN

# Decision Tree

## Confusion Matrix for Decision Tree

| True label \ Predicted label | will not change | will change |
|---|---|---|
| will not change | 987 | 134 |
| will change | 232 | 147 |

# Logistic Regression:

## Confusion Matrix for Logistic Regression

| True label \ Predicted label | will not change | will change |
|---|---|---|
| will not change | 987 | 134 |
| will change | 232 | 147 |

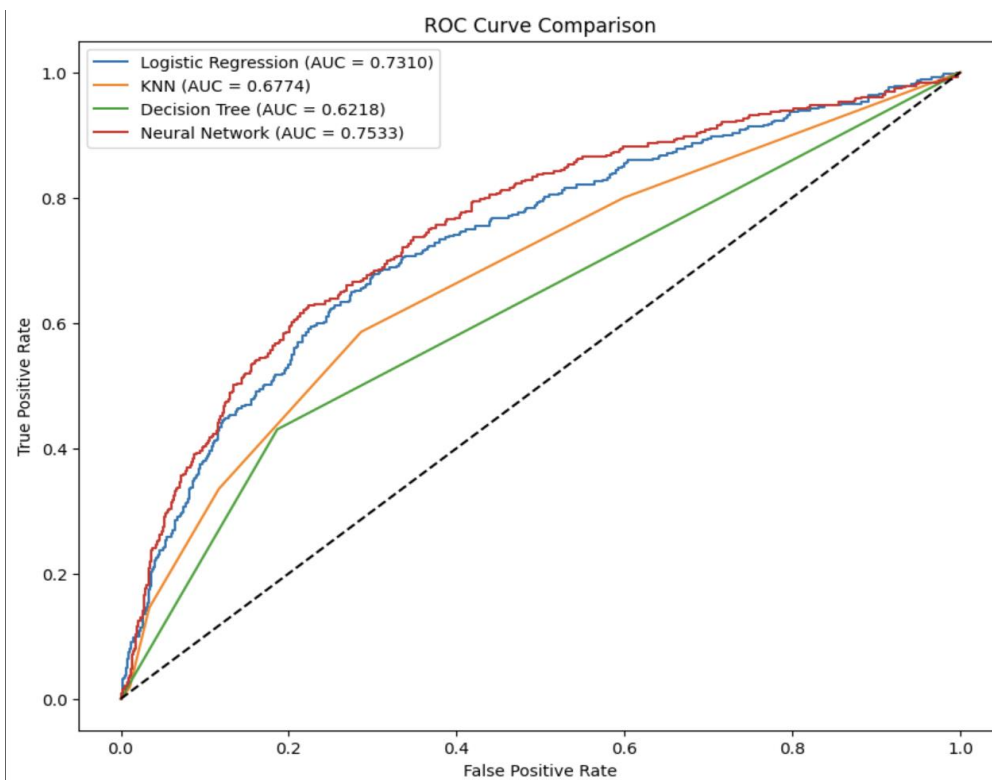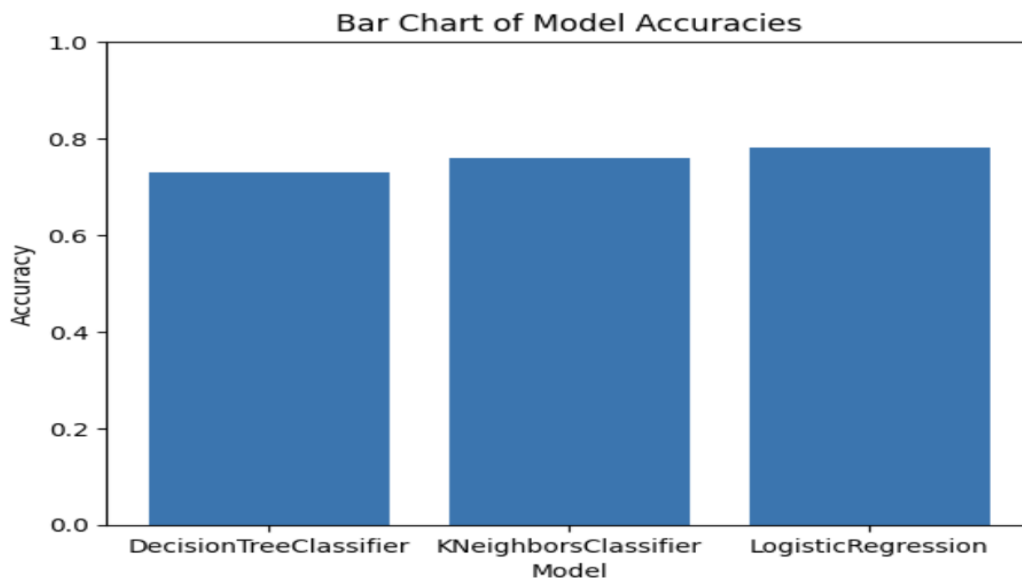## Loss Curve of Neural Network



## AUC score, ROC curve

**Accuracy Bar Chart:**



**8. Conclusion:** After evaluating three different machine learning models — Decision Tree, K-Nearest Neighbors (KNN), and Logistic Regression — on the career switch prediction classification problem, the following accuracies were obtained: Decision Tree: 72%, K-Nearest Neighbors (KNN): 74%, and Logistic Regression: 77% and Neural Networks:                                                                                   78%.
 Among these models, Logistic Regression achieved the highest accuracy and proved to be the most effective for this prediction task. However, these results are subject to change with different datasets, and further improvement is possible through hyperparameter tuning, feature engineering, and testing with additional algorithms.

# Reference:

TensorFlow. (n.d.). *Getting started with Keras & Sequential model.* TensorFlow.
https://www.tensorflow.org/guide/keras/sequential_model

Scikit-learn developers. (n.d.). *User guide: Classification*. Scikit-learn.
https://scikit-learn.org/stable/supervised_learning.html

Towards Data Science. (2020, July 20). *Understanding ROC Curves and AUC in Machine Learning.* Medium.
https://towardsdatascience.com/understanding-roc-curves-and-auc-840b6c2839b7

GeeksforGeeks. (n.d.). *Logistic Regression in Machine Learning*. GeeksforGeeks.
https://www.geeksforgeeks.org/logistic-regression/

GeeksforGeeks. (n.d.). *K-Nearest Neighbors (KNN) Algorithm in Python.* GeeksforGeeks. https://www.geeksforgeeks.org/k-nearest-neighbours/